



# Differentially Private Chi-Squared Hypothesis Testing: Goodness of Fit and Independence Testing

## Citation

Gaboardi, Marco, Hyun-Woo Lim, Ryan M. Rogers, and Salil P. Vadhan. 2016. "Differentially Private Chi-Squared Hypothesis Testing: Goodness of Fit and Independence Testing." In ICML'16 Proceedings of the 33rd International Conference on International Conference on Machine Learning, New York, NY, June 19-24, 2016, Volume 48: 2111-2120.

## Permanent link

<http://nrs.harvard.edu/urn-3:HUL.InstRepos:34614371>

## Terms of Use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Open Access Policy Articles, as set forth at <http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP>

## Share Your Story

The Harvard community has made this article openly available.  
Please share how this access benefits you. [Submit a story](#).

[Accessibility](#)

# Differentially Private Chi-Squared Hypothesis Testing: Goodness of Fit and Independence Testing\*

Marco Gaboardi<sup>†1</sup>, Hyun woo Lim<sup>2</sup>, Ryan Rogers<sup>3</sup>, and Salil P. Vadhan<sup>‡4</sup>

<sup>1</sup>*University at Buffalo, SUNY*

<sup>2</sup>*University of California, Los Angeles*

<sup>3</sup>*University of Pennsylvania*

<sup>4</sup>*Harvard University*

May 25, 2016

## Abstract

Hypothesis testing is a useful statistical tool in determining whether a given model should be rejected based on a sample from the population. Sample data may contain sensitive information about individuals, such as medical information. Thus it is important to design statistical tests that guarantee the privacy of subjects in the data. In this work, we study hypothesis testing subject to differential privacy, specifically chi-squared tests for goodness of fit for multinomial data and independence between two categorical variables.

We propose new tests for goodness of fit and independence testing that like the classical versions can be used to determine whether a given model should be rejected or not, and that additionally can ensure differential privacy. We give both Monte Carlo based hypothesis tests as well as hypothesis tests that more closely follow the classical chi-squared goodness of fit test and the Pearson chi-squared test for independence. Crucially, our tests account for the distribution of the noise that is injected to ensure privacy in determining significance.

We show that these tests can be used to achieve desired significance levels, in sharp contrast to direct applications of classical tests to differentially private contingency tables which can result in wildly varying significance levels. Moreover, we study the statistical power of these tests. We empirically show that to achieve the same level of power as the classical non-private tests our new tests need only a relatively modest increase in sample size.

---

\*This work is part of the “Privacy Tools for Sharing Research Data” project based at Harvard, supported by NSF grant CNS-1237235 as well as a grant from the Sloan Foundation.

<sup>†</sup>This work has been partially supported by the “PrivInfer - Programming Languages for Differential Privacy: Conditioning and Inference” EPSRC project EP/M022358/1 and by the University of Dundee, UK.

<sup>‡</sup>Also supported by a Simons Investigator grant. Work done in part while visiting the Department of Applied Mathematics and the Shing-Tung Yau Center at National Chiao-Tung University in Taiwan.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Contributions . . . . .	3
1.2	Hypothesis testing for the social sciences . . . . .	4
<b>2</b>	<b>Related Work</b>	<b>4</b>
<b>3</b>	<b>Differential Privacy Preliminaries</b>	<b>6</b>
<b>4</b>	<b>Hypothesis Testing Preliminaries</b>	<b>7</b>
<b>5</b>	<b>Goodness of Fit Test</b>	<b>7</b>
5.1	Differentially Private Chi-Squared Statistic . . . . .	9
5.2	Monte Carlo Test: $\text{MCGOF}_{\mathcal{D}}$ . . . . .	10
5.3	Asymptotic Approach: $\text{PrivGOF}$ . . . . .	11
5.4	Power Analysis of $\text{PrivGOF}$ . . . . .	12
<b>6</b>	<b>Independence Testing</b>	<b>14</b>
6.1	Estimating Parameters with Private Counts . . . . .	17
6.2	Monte Carlo Test: $\text{MCIndep}_{\mathcal{D}}$ . . . . .	18
6.3	Asymptotic Approach: $\text{PrivIndep}$ . . . . .	19
<b>7</b>	<b>Significance Results</b>	<b>19</b>
<b>8</b>	<b>Power Results</b>	<b>21</b>
<b>9</b>	<b>Conclusion</b>	<b>22</b>
	<b>References</b>	<b>24</b>

# 1 Introduction

Hypothesis testing provides a systematic way to test given models based on a sample, so that with high confidence a data analyst may conclude that the model is incorrect or not. However, these data samples may contain highly sensitive information about the subjects and so the privacy of individuals can be compromised when the results of a data analysis are released. For example, in the area of *genome-wide association studies* (GWAS) Homer et al. [2008] have shown that it is possible to identify subjects in a data set based on publicly available aggregate statistics.

A way to address this concern is by developing new techniques to support privacy-preserving data analysis. An approach that is gaining more and more attention by the statistics and data analysis community is *differential privacy* [Dwork et al., 2006b], which originated in theoretical computer science. In this work, we seek to develop hypothesis tests that are differentially private and that give conclusions similar to standard, non-private hypothesis tests.

We focus here on two classical tests for data drawn from a multinomial distribution: *goodness of fit test*, which determines whether the data was in fact drawn from a multinomial distribution with probability vector  $\mathbf{p}^0$ ; and *independence test*, which tests whether two categorical random variables are independent of each other. Both tests depend on the *chi-squared* statistic, which is used to determine whether the data is likely or not under the given model.

To guarantee differential privacy, we consider adding Laplace and Gaussian noise to the counts of categorical data. Using the noisy data we can form a *private* chi-squared statistic. It turns out that the classical hypothesis tests perform poorly when used with this modified statistic because they ignore the fact that noise was added. To improve this situation, we develop four new tests that account for the additional noise due to privacy. In particular, we give two tests based on a *Monte Carlo* approach to testing the null hypothesis and two tests based on an *asymptotic* distribution of the private chi-squared distribution factoring in the noise distribution.

Our four differentially private tests achieve a target level  $1 - \alpha$  *significance*, i.e. they reject with probability at most  $\alpha$  when the null hypothesis holds (in some cases, we provide a rigorous proof of this fact and in others, it is experimentally verified). This guarantees limited Type I errors. However, all of our tests do lose *power*; that is when the null hypothesis is false, they correctly reject with lower probability than the classical hypothesis tests. This corresponds to an increase of Type II errors. We empirically show that we can recover a level of power similar to the one achieved by the classical versions by adding more samples.

## 1.1 Contributions

For goodness of fit testing we present two differentially private tests. First, we give a Monte Carlo (MC) based test  $\text{MCGOF}_{\mathcal{D}}$  that guarantees significance at least  $1 - \alpha$  for any desired  $\alpha > 0$  (commonly 0.05 or 0.10), when either Laplace or Gaussian noise is added in the private chi-squared statistic. When Gaussian noise is used in the private chi-squared statistic and has privacy parameter that decays with the sample size, we then analytically obtain an asymptotic distribution for the private chi-squared statistic that is a linear combination of independent chi-squared random variables with one degree of freedom, given that the null hypothesis is true. We then use this asymptotic distribution for the private chi-squared statistic for the test  $\text{PrivGOF}$ . This provides an alternative to our MC test, which is computationally less expensive and more closely parallels the classical test, which is based on the asymptotic distribution of the nonprivate statistic (which has an asymptotic chi-squared distribution with  $d - 1$  degrees of freedom, where  $d$  is the dimension of the data).

Further, when the data actually satisfies an alternate hypothesis  $\mathbf{p}^1 \neq \mathbf{p}^0$  that has a specific form, we find the asymptotic distribution of the private chi-squared, which is to be compared with the non-private chi-squared statistic converging to a non-central chi-squared distribution with  $d - 1$  degrees of freedom when the alternate hypothesis is true.

We then turn to independence testing. Given a contingency table where each cell has Laplace or Gaussian noise added, inspired by Karwa and Slavković [2016], we present a (heuristic) procedure 2MLE for finding an approximate maximum likelihood estimator (MLE) for the true probability vector that satisfies the independence hypothesis. Note that for the classical Pearson chi-squared test, one computes the MLE for the true probability vector given a contingency table (without additional noise). We then use this estimate as the probability vector in our private chi-squared statistic and give an MC method **MCIndep<sub>D</sub>** based on the private chi-squared statistic.

Lastly, when dealing with Gaussian noise, we give a differentially private test **PrivIndep** that closely follows the analysis of the Pearson chi-squared test for independence. We show that when we approximate the distribution of the private chi-squared statistic with a linear combination of chi-squared random variables with 1 degree of freedom we obtain empirical significance at least  $1 - \alpha$  and power that closely follows the power of **MCIndep<sub>D</sub>** that uses Gaussian noise for various sample sizes  $n$ .

For all of our tests we give empirical significance and power results and compare them to their corresponding classical, non-private tests. We obtain empirical significances that are near the desired  $1 - \alpha$  level and clearly outperform the classical tests when they are used on the noisy data. In particular, this guarantees that our tests do not incur more Type I error than the classical tests. However, our tests do have lower power than the corresponding non-private test, due to the additional noise injected. Thus, our tests have larger Type II errors than the classical (nonprivate) tests. In particular, the tests using Gaussian noise have a larger loss in power than the ones using Laplace noise. This is because for a given level of privacy, the Gaussian noise has a larger variance than the Laplace noise. To achieve power similar to the classical tests, our private tests require more samples.

## 1.2 Hypothesis testing for the social sciences

Our work is part of the broader effort of the project “Privacy Tools for Sharing Research Data”<sup>1</sup> that aims in particular at developing differentially private tools that can be used for studies in the social sciences. Social scientists often deal with various sensitive data that contains individual’s private information, e.g. voting behavior [Greenwald et al., 1987], attitude toward abortion [Ebaugh and Haney, 1978] and medical records [David and Beards, 1985]. The framework of hypothesis testing is frequently used by social scientists to confirm or reject their belief to how a population is modeled, e.g. goodness of fit tests have been used by [David and Beards, 1985, Gill et al., 1987, Blair et al., 1979, Glaser, 1959] and independence tests have been used by [Kuklinski and West, 1981, Ebaugh and Haney, 1978, Berry, 1961, Krain and Myers, 1997, Greenwald et al., 1987, Mitchell and McCormick, 1988].

---

<sup>1</sup><http://privacytools.seas.harvard.edu>

## 2 Related Work

There has been a myriad of work dealing with the application of differential privacy in statistical inference. One of the first works that put differential privacy in the language of statistics is Wasserman and Zhou [2010], which studies rates of convergence of distributions based on differentially private data released from the *exponential mechanism* [McSherry and Talwar, 2007]. In a result of great generality, Smith [2011] shows that for a wide class of statistics  $T$ , there is a differentially private statistic that converges in distribution to the same asymptotic distribution as  $T$ . However, having the correct asymptotic distribution does not ensure that only statistically significant conclusions are drawn at finite sample sizes, and indeed we observe that this fails dramatically for the most natural differentially private algorithms. Thus, we study how to ensure significance and optimize power at small sample sizes by focusing on two basic statistical tests.

A tempting first approach to developing a hypothesis test for categorical data that is also differentially private is to either add noise directly to the chi-squared statistic that will ensure differential privacy or to add noise to each cell count (as we do in this work) and use a classical test with the private counts. For the former method, the amount of noise that must be added to ensure privacy can be unbounded in the worst case. However, motivated by applications to genome-wide association studies (GWAS), Uhler et al. [2013] and Yu et al. [2014] place restrictions on the form of the data or what is known to the data analyst to reduce the scale of the noise that needs to be added. The work of Johnson and Shmatikov [2013] adds noise to each cell of a contingency table, but then uses classical statistical tests on the private version of the data, which we show can have very poor significance (see Figures 1 and 2). Additionally, Uhler et al. [2013] look at  $3 \times 2$  contingency tables that are evenly split between the two columns, and study releasing differentially private  $\chi^2$ -statistics of the most relevant SNPs for certain diseases by perturbing the table of counts, the  $\chi^2$ -statistic itself, and the  $p$ -values for the underlying test. The only one of these works that explicitly examine significance and power in hypothesis testing (as we do here) is Uhler et al. [2013], which shows that perturbing the  $p$ -values in independence testing does not perform much better than a random test, independent of a selected threshold, e.g.  $\alpha$ . In fact, Uhler et al. [2013] goes as far as to say that basing inference on perturbed  $p$ -values “seems impossible.” An interesting direction for future work would be to apply the *distance-score mechanism* introduced by Johnson and Shmatikov [2013] and later improved by Yu et al. [2014], Yu and Ji [2014], Simmons and Berger [2016], to achieving a target level of significance and high power in hypothesis testing for GWAS data.

If we assume that there is some prior estimates for the contingency table cell probabilities, Vu and Slavkovic [2009] determine the sample size adjustment for the Pearson chi-squared independence test that uses the private counts to achieve the same power as the test with the original counts. Several other works have shown negative experimental results for using classical inference on statistics that have been altered for differential privacy [Fienberg et al., 2010, Karwa and Slavković, 2012, Karwa and Slavković, 2016].

Another problem that arises when noise is added to the cells in a contingency table is that the entries may neither be positive nor sum to a known value  $n$ . Several works have focused on this problem, where they seek to release a contingency table in a differentially private way that also satisfies some known *consistency* properties of the underlying data [Barak et al., 2007, Li et al., 2010, Hardt et al., 2012, Li and Miklau, 2012, Gaboardi et al., 2014]. For independence testing, we use techniques from Lee et al. [2015] to find the most likely contingency table given the noisy version of it so that we can then estimate the cell probabilities that generated the table. This two

step procedure to estimate parameters given a differentially private statistic is inspired by the work of Karwa and Slavković [2016] for estimating parameters in the  $\beta$ -model for random graphs.

Independent of our work, Wang et al. [2015] also look at hypothesis testing with categorical data subject to differential privacy. They mainly consider adding Laplace noise to the data but point out that their method also generalizes to arbitrary noise distributions. However, in order to compute critical values, they resort to Monte Carlo methods to sample from the asymptotic distribution. Our Monte Carlo approach samples from the *exact* distribution from the underlying null hypothesis, which, unlike sampling from the asymptotic distribution, guarantees significance at least  $1 - \alpha$  in goodness of fit tests at finite sample sizes. We only focus on Gaussian noise in our asymptotic analysis due to there being existing methods for finding tail probabilities (and hence critical values) for the resulting distributions, but our approaches can be generalized for arbitrary noise distributions. Further, we consider the power of each of our differentially private tests.

### 3 Differential Privacy Preliminaries

We start with a brief overview of differential privacy. In order to define differential privacy, we first define *neighboring* databases  $\mathbf{d}, \mathbf{d}'$  from some class of databases  $D^n$  where they differ in an individual's data but are equal among the rest of the data, e.g.  $\mathbf{d} = (d_1, \dots, d_i, \dots, d_n)$  and  $\mathbf{d}' = (d_1, \dots, d'_i, \dots, d_n)$  where  $d_i \neq d'_i$ . We will consider  $n$  to be known and public.

**Definition 3.1** (Differential Privacy [Dwork et al., 2006b]). Let  $M : D^n \rightarrow O$  be some randomized mechanism. For  $\epsilon, \delta > 0$  we say that  $M$  is  $(\epsilon, \delta)$ -differentially private if for any neighboring databases  $\mathbf{d}, \mathbf{d}' \in D^n$  and any subset of outcomes  $S \subseteq O$  we have

$$\Pr[M(\mathbf{d}) \in S] \leq e^\epsilon \Pr[M(\mathbf{d}') \in S] + \delta.$$

If  $\delta = 0$ , then we simply say  $M$  is  $\epsilon$ -differentially private. The meaning of these parameters, loosely, is that with probability  $1 - \delta$  there is at most  $\epsilon$  information leakage (so with probability at most  $\delta$  it can leak lots of information). For this reason, we will think of  $\epsilon$  as a small constant, e.g. 0.1, and  $\delta \ll 1/n$  as cryptographically small, where we sometimes write  $\delta_n$  to explicitly show its dependence on  $n$ .

A typical differentially private mechanism is to add carefully calibrated noise to some quantity that a data analyst is interested in. We can release a differentially private answer to a function  $\phi : D^n \rightarrow \mathbb{R}^d$  by adding independent noise to each component of  $\phi$ . The scale of the noise we add depends on the impact any individual can have on the outcome. We use the *global sensitivity* of  $\phi$  to quantify this impact, which we define for  $i = 1, 2$  as:

$$GS_i(\phi) = \max_{\mathbf{d}, \mathbf{d}' \text{ neighboring in } D^n} \{ \|\phi(\mathbf{d}) - \phi(\mathbf{d}')\|_i \}.$$

**Lemma 3.2** (Dwork et al. [2006a,b]). Let  $\phi : D^n \rightarrow \mathbb{R}^d$  have global sensitivity  $GS_i(\phi)$  for  $i = 1, 2$ . Then the mechanism  $M_{\mathcal{D}} : D^n \rightarrow \mathbb{R}^d$  where

$$M_{\mathcal{D}}(\mathbf{d}) = \phi(\mathbf{d}) + (Z_1, \dots, Z_d)^T \quad \{Z_i\} \stackrel{i.i.d.}{\sim} \mathcal{D}$$

is  $\epsilon$ -differentially private if  $\mathcal{D} = \text{Laplace}\left(\frac{GS_1(\phi)}{\epsilon}\right)$  or  $(\epsilon, \delta)$ -differentially private if  $\mathcal{D} = N(0, \sigma^2)$  with  $\sigma = \frac{GS_2(\phi)\sqrt{2\ln(2/\delta)}}{\epsilon}$ .

There are many useful properties of differentially private mechanisms. The one we will use in this paper is referred to as *post-processing*, which ensures privacy no matter what we do with the outcome of  $M$ .

**Lemma 3.3.** [*Post-Processing [Dwork et al., 2006b]*] Let  $M : D^n \rightarrow O$  be  $(\epsilon, \delta)$ -differentially private and  $\psi : O \rightarrow O'$  be some arbitrary mapping from  $O$  to  $O'$ . Then  $\psi \circ M : D^n \rightarrow O'$  remains  $(\epsilon, \delta)$ -differentially private.

The tests that we present will be differentially private, assuming  $n$  is known and public, because we will add Laplace or Gaussian noise as in Lemma 3.2 to the vector of counts in goodness of fit testing

## 4 Hypothesis Testing Preliminaries

Given sampled data from a population, we wish to test whether the data came from a specific model, which is given as a *null hypothesis*  $H_0$ . We will denote our test as an algorithm  $\mathcal{A}$  that takes a dataset  $\mathbf{X}$ , significance level  $1 - \alpha$  and null hypothesis  $H_0$  and returns a decision of whether to reject  $H_0$  or not. We would like to design our test so that we achieve Type I error at most  $\alpha$ , that is

$$\Pr [\mathcal{A}(\mathbf{X}; \alpha, H_0) = \text{Reject} | H_0] \leq \alpha$$

while also achieving a small Type II error  $\beta = \Pr [\mathcal{A}(\mathbf{X}; \alpha, H_0) = \text{Reject} | H_1]$  when the model is actually some *alternate*  $H_1 \neq H_0$ . We think of bounding Type I error as a *hard constraint* in our tests and then hope to minimize Type II error. Note that the probability is taken over the randomness from the data generation and the possible randomness from the algorithm  $\mathcal{A}$  itself. It is common to refer to  $1 - \alpha$  as the *significance* of test  $\mathcal{A}$  and  $1 - \beta$  as the *power* of  $\mathcal{A}$ .

## 5 Goodness of Fit Test

We consider  $\mathbf{X} = (X_1, \dots, X_d)^T \sim \text{Multinomial}(n, \mathbf{p})$  where  $\mathbf{p} = (p_1, \dots, p_d)$  and  $\sum_{i=1}^d p_i = 1$ . Note that the multinomial distribution is the generalization of a binomial distribution where there are  $d$  outcomes. For a *goodness of fit test*, we want to test the null hypothesis  $H_0 : \mathbf{p} = \mathbf{p}^0$ . A common way to test this is based on the *chi-squared* statistic  $Q^2$  where

$$Q^2 = \sum_{i=1}^d \frac{(X_i - np_i^0)^2}{np_i^0} \quad (1)$$

We present the classical chi-squared *goodness of fit test* in Algorithm 1, which compares the chi-squared statistic  $Q^2$  to a threshold  $\chi_{d-1, 1-\alpha}^2$  that depends on a desired level of significance  $1 - \alpha$  as well as the dimension of the data. The threshold  $\chi_{d-1, 1-\alpha}^2$  satisfies the following relationship:

$$\Pr [\chi_{d-1}^2 \geq \chi_{d-1, 1-\alpha}^2] = \alpha.$$

where  $\chi_{d-1}^2$  is a chi-squared random variable with  $d-1$  degrees of freedom, which is the distribution of the random variable  $\mathbf{N}^T \mathbf{N}$  where  $\mathbf{N} \sim N(0, I_{d-1})$ .

The reason why we compare  $Q^2$  with the chi-squared distribution is because of the following classical result.



---

**Algorithm 1** Goodness of Fit Test for Multinomial Data

---

```

procedure GOF(Data  $\mathbf{x}$ , Significance  $1 - \alpha$ , and  $H_0 : \mathbf{p} = \mathbf{p}^0$ )
  Compute  $Q^2$  from (1)
  if  $Q^2 > \chi_{d-1, 1-\alpha}^2$  then
    Decision  $\leftarrow$  Reject
  else
    Decision  $\leftarrow$  Fail to Reject
  return Decision.

```

---

**Theorem 5.1.** [Bishop et al., 1975] Assuming  $H_0 : \mathbf{p} = \mathbf{p}^0$  holds, the statistic  $Q^2$  converges in distribution to a chi-squared with  $d - 1$  degrees of freedom, i.e.

$$Q^2 \xrightarrow{D} \chi_{d-1}^2.$$

Note that this does not guarantee that  $\Pr [Q^2 > \chi_{d-1, 1-\alpha}^2] \leq \alpha$  for finite samples, nevertheless the test works well and is widely used in practice.

It will be useful for our purposes to understand why the asymptotic result holds in Theorem 5.1. We present the following classical analysis [Bishop et al., 1975] of Theorem 5.1 so that we can understand what adjustments need to be made to find an approximate distribution for a differentially private statistic. Consider the random vector  $\mathbf{U} = (U_1, \dots, U_d)$  where

$$U_i = \frac{X_i - np_i^0}{\sqrt{np_i^0}} \quad \forall i \in [d]. \quad (2)$$

We write the covariance matrix for  $\mathbf{U}$  as  $\Sigma$  where

$$\Sigma = I_d - \sqrt{\mathbf{p}^0} \sqrt{\mathbf{p}^0}^T \quad (3)$$

and  $\sqrt{\mathbf{p}^0} = (\sqrt{p_1^0}, \dots, \sqrt{p_d^0})^T$ . By the *central limit theorem* we know that  $\mathbf{U}$  converges in distribution to a multivariate normal

$$\mathbf{U} \xrightarrow{D} N(\mathbf{0}, \Sigma) \quad \text{as } n \rightarrow \infty.$$

Thus, when we make the assumption that  $\mathbf{U}$  is multivariate normal, then the significance of GOF given in Algorithm 1 is exactly  $1 - \alpha$ .

We show in the following lemma that if a random vector  $\mathbf{U}$  is exactly distributed as multivariate normal then we get that  $Q^2 = \mathbf{U}^T \mathbf{U} \sim \chi_{d-1}^2$ .

**Lemma 5.2** ([Bishop et al., 1975]). *If  $\mathbf{U} \sim N(\mathbf{0}, \Sigma)$  for  $\Sigma$  given in (3) then  $\mathbf{U}^T \mathbf{U} \sim \chi_{d-1}^2$ .*

*Proof.* The eigenvalues of  $\Sigma$  must be either 0 or 1 because  $\Sigma$  is idempotent. Thus, the number of eigenvalues that are 1 equals trace of  $\Sigma$ , which is  $d - 1$ . We then know that there exists a matrix  $H \in \mathbb{R}^{d \times d-1}$  where  $\Sigma = HH^T$  and  $H^T H = I_{d-1}$ . Define the random variable  $\mathbf{Y} \sim N(\mathbf{0}, I_{d-1})$ . Note that  $H\mathbf{Y}$  is equal in distribution to  $\mathbf{U}$ . We then have

$$\mathbf{U}^T \mathbf{U} \sim \mathbf{Y}^T H^T H \mathbf{Y} \sim \mathbf{Y}^T \mathbf{Y} \sim \chi_{d-1}^2$$

□

## 5.1 Differentially Private Chi-Squared Statistic

To ensure differential privacy, we add independent noise to each component of  $\mathbf{X}$ , which we will either use Laplace or Gaussian noise. The function  $g$  that outputs the counts in the  $d$  cells has global sensitivity  $GS_1(g) = 2$  and  $GS_2(g) = \sqrt{2}$  because one individual may move from one cell count (decreasing the count by 1) to another (increasing the cell count by 1). We then form the *private chi-squared* statistic  $Q_{\mathcal{D}}^2$  based on the noisy counts,

$$Q_{\mathcal{D}}^2 = \sum_{i=1}^d \frac{(X_i + Z_i - np_i^0)^2}{np_i^0}, \quad \{Z_i\} \stackrel{i.i.d.}{\sim} \mathcal{D} \quad (4)$$

where the distributions for the noise that we consider include

$$\mathcal{D} = \text{Laplace}(2/\epsilon) \quad \text{and} \quad \mathcal{D} = N(0, \sigma^2) \quad \text{where} \quad \sigma(\epsilon, \delta_n) = \frac{2\sqrt{\ln(2/\delta_n)}}{\epsilon}. \quad (5)$$

We will denote the  $\epsilon$ -differentially private statistic as  $Q_{\text{Lap}}^2$  and the  $(\epsilon, \delta)$ -differentially private statistic as  $Q_{\text{Gauss}}^2$  based on whether we use Laplace or Gaussian noise, respectively. Recall that in the original goodness of fit test without privacy in Algorithm 1 we compare the distribution of  $Q^2$  with that of a chi-squared random variable with  $d - 1$  degrees of freedom. The following result shows that adding noise to each cell count does not affect this asymptotic distribution.

**Lemma 5.3.** *Fixing  $\mathbf{p}^0 > 0$ , and having privacy parameters  $(\epsilon, \delta_n)$  where  $\epsilon > 0$  and  $\delta_n$  satisfies the following condition  $\frac{\log(2/\delta_n)}{n\epsilon^2} \rightarrow 0$ , then the private chi-squared statistic  $Q_{\mathcal{D}}^2$  given in (4) converges in distribution to  $\chi_{d-1}^2$  as  $n \rightarrow \infty$ .*

*Proof.* We first expand (4) to get

$$Q_{\mathcal{D}}^2 = \sum_{i=1}^d \left( \frac{X_i - np_i^0}{\sqrt{np_i^0}} \right)^2 + 2 \sum_{i=1}^d \left( \frac{Z_i}{\sqrt{np_i^0}} \right) \left( \frac{X_i - np_i^0}{\sqrt{np_i^0}} \right) + \sum_{i=1}^d \left( \frac{Z_i}{\sqrt{np_i^0}} \right)^2$$

We first focus on  $\mathcal{D}$  being Gaussian. We define the two random vectors  $\mathbf{Z}^{(n)} = \left( \frac{Z_i}{\sqrt{np_i^0}} \right)_{i=1}^n$  and  $\mathbf{X}^{(n)} = \left( \frac{X_i - np_i^0}{\sqrt{np_i^0}} \right)_{i=1}^n$ . We have that  $\text{Var}(Z_i^{(n)}) = \frac{\sigma(\epsilon, \delta_n)^2}{np_i^0} = \frac{2\ln(2/\delta_n)}{np_i^0 \epsilon^2}$  which goes to zero by hypothesis. Additionally  $\mathbb{E}[\mathbf{Z}^{(n)}] = 0$ , so we know that  $\mathbf{Z}^{(n)} \xrightarrow{P} \mathbf{0}$ . We also know that  $\mathbf{X}^{(n)} \xrightarrow{D} N(0, \Sigma)$ , so that  $\mathbf{Z}^{(n)} \cdot \mathbf{X}^{(n)} \xrightarrow{D} 0$  by Slutsky's Theorem<sup>2</sup> and thus  $\mathbf{Z}^{(n)} \cdot \mathbf{X}^{(n)} \xrightarrow{P} 0$  (because 0 is constant). Another application of Slutsky's Theorem tells us that  $Q_{\mathcal{D}}^2 \xrightarrow{D} \chi_{d-1}^2$ , since  $Q_{\mathcal{D}}^2 - Q^2 \xrightarrow{P} 0$  and  $Q^2 \xrightarrow{D} \chi_{d-1}^2$ . The proof for  $\mathcal{D}$  being Laplacian follows the same analysis.  $\square$

It then seems natural to use GOF on the private chi-squared statistic as if we had the actual chi-squared statistic that did not introduce noise to each count since both private and nonprivate statistics have the same asymptotic distribution. We will show in our results in Section 7 that if we were to simply compare the private statistic to the critical value  $\chi_{d-1, 1-\alpha}^2$ , we will typically not get

<sup>2</sup>Slutsky's Theorem states that if  $X_n \xrightarrow{D} X$  and  $Z_n \xrightarrow{P} c$  then  $X_n \cdot Z_n \xrightarrow{D} cX$  and  $X_n + Z_n \xrightarrow{D} X + c$ .

a good significance level even for relatively large  $n$  which we need in order for it to be practical tool for data analysts. In the following lemma we show that for every realization of data, the statistic  $Q_{\mathcal{D}}^2$  is expected to be larger than the actual chi-squared statistic  $Q^2$ .

**Lemma 5.4.** *For each realization  $\mathbf{X} = \mathbf{x}$ , we have  $\mathbb{E}_{\mathcal{D}} [Q_{\mathcal{D}}^2 | \mathbf{x}] \geq Q^2$ , where  $\mathcal{D}$  has mean zero.*

*Proof.* Consider the convex function  $f(y) = y^2$ . Applying Jensen's inequality, we have  $f(y) \leq \mathbb{E}_{Z_i} [f(y + Z_i)]$  for all  $i = 1, \dots, d$  where  $Z_i$  is sampled i.i.d. from  $\mathcal{D}$  which has mean zero. We then have for  $\mathbf{X} = \mathbf{x}$

$$\begin{aligned} Q^2 &= \sum_{i=1}^d \frac{f(x_i - np_i^0)}{np_i^0} = \sum_{i=1}^d \frac{f(\mathbb{E}[x_i - np_i^0 + Z_i])}{np_i^0} \\ &\leq \mathbb{E}_{\{Z_i\} \stackrel{i.i.d.}{\sim} \mathcal{D}} \left[ \sum_{i=1}^d \frac{f(x_i - np_i^0 + Z_i)}{np_i^0} \right] = \mathbb{E}_{\{Z_i\} \stackrel{i.i.d.}{\sim} \mathcal{D}} [Q_{\mathcal{D}}^2 | \mathbf{x}] \end{aligned}$$

□

This result suggests that the significance threshold for the private version of the chi-squared statistic  $Q_{\mathcal{D}}^2$  should be higher than the standard one. Otherwise, we would reject  $H_0$  too easily using the classical test, which we show in our experimental results. This motivates the need to develop new tests that account for the distribution of the noise.

## 5.2 Monte Carlo Test: MCGOF $_{\mathcal{D}}$

Given some null hypothesis  $\mathbf{p}^0$  and statistic  $Q_{\mathcal{D}}^2$ , we want to determine a threshold  $\tau^\alpha$  such that  $Q_{\mathcal{D}}^2 > \tau^\alpha$  at most an  $\alpha$  fraction of the time when the null hypothesis is true. As a first approach, we determine threshold  $\tau^\alpha$  using a Monte Carlo (MC) approach by sampling from the distribution of  $Q_{\mathcal{D}}^2$ , where  $\mathbf{X} \sim \text{Multinomial}(n, \mathbf{p}^0)$  and  $\{Z_i\} \stackrel{i.i.d.}{\sim} \mathcal{D}$  for both Laplace and Gaussian noise.

Let  $M_1, \dots, M_k$  be  $k$  continuous random variables that are i.i.d. from the distribution of  $Q_{\mathcal{D}}^2$  assuming  $H_0$ . Further let  $M$  be a fresh sample from the distribution of  $Q_{\mathcal{D}}^2$  assuming  $H_0$ . We will write the density and distribution of  $Q_{\mathcal{D}}^2$  as  $f(\cdot)$  and  $F(\cdot)$ , respectively. Our test will reject  $M$  if it falls above some threshold, i.e. critical value, which we will take to be the  $t$ -th order statistic of  $\{M_i\}$ , also written as  $M_{(t)}$ , so that with probability at most  $\alpha$ ,  $M$  is above this threshold. This will guarantee significance at least  $1 - \alpha$ . We then find the smallest  $t \in [k]$  such that  $\alpha \geq \Pr [M > M_{(t)}]$ , or

$$\begin{aligned} \alpha &\geq \int_{-\infty}^{\infty} f(m) \sum_{j=t}^k \binom{k}{j} F(m)^j (1 - F(m))^{k-j} dm = \int_0^1 \sum_{j=t}^k \binom{k}{j} p^j (1 - p)^{k-j} dp \\ &= \sum_{j=t}^k \frac{1}{k+1} \implies t \geq (k+1)(1 - \alpha). \end{aligned}$$

We then set our threshold based on the  $\lceil (k+1)(1 - \alpha) \rceil$  ordered statistic of our  $k$  samples. By construction, this will ensure that we achieve the significance level we want. Our test then is to sample  $k$  points from the distribution of  $Q_{\mathcal{D}}^2$  and then take the  $\lceil (k+1)(1 - \alpha) \rceil$ -percentile as our cutoff, i.e. if our statistic falls above this value, then we reject  $H_0$ . Note that we require  $k \geq 1/\alpha$ , otherwise there would not be a  $\lceil (k+1)(1 - \alpha) \rceil$  ordered statistic in  $k$  samples. We give the resulting test in Algorithm 2.

---

**Algorithm 2** MC Goodness of Fit

---

```
procedure MCGOF $_{\mathcal{D}}$ (Data  $\mathbf{x}$ ; Privacy  $(\epsilon, \delta)$ , Significance  $1 - \alpha$ ,  $H_0 : \mathbf{p} = \mathbf{p}^0$ )  
  Compute  $q = Q_{\mathcal{D}}^2(4)$ .  
  Select  $k > 1/\alpha$ .  
  Sample  $k$  points  $q_1, \dots, q_k$  i.i.d. from the distribution of  $Q_{\mathcal{D}}^2$  and sort them  $q_{(1)} \leq \dots \leq q_{(k)}$ .  
  Compute threshold  $q_{(t)}$  where  $t = \lceil (k+1)(1-\alpha) \rceil$ .  
  if  $q > q_{(t)}$  then  
    Decision  $\leftarrow$  Reject  
  else  
    Decision  $\leftarrow$  Fail to Reject  
  return Decision.
```

---

**Theorem 5.5.** *The test  $MCGOF_{\mathcal{D}}(\cdot, (\epsilon, \delta), \alpha, \mathbf{p}^0)$  has significance at least  $1 - \alpha$ , also written as  $\Pr [MCGOF_{\mathcal{D}}(\mathbf{X}, (\epsilon, \delta), \alpha, \mathbf{p}^0) = \text{Reject} | H_0] \leq \alpha$ .*

In Section 7, we present the empirical power results for  $MCGOF_{\mathcal{D}}$  (along with all our other tests) when we fix an alternative hypothesis.

### 5.3 Asymptotic Approach: PrivGOF

In this section we attempt to determine a more analytical approximation to the distribution of  $Q_{Gauss}^2$ . We focus on Gaussian noise because it is more compatible with the asymptotic analysis of GOF, as opposed to Laplace noise. Recall the random vector  $\mathbf{U}$  given in (2). We then introduce the Gaussian noise random vector as  $\mathbf{V} = (Z_1/\sigma(\epsilon, \delta_n), \dots, Z_d/\sigma(\epsilon, \delta_n))^T \sim N(\mathbf{0}, I_d)$ . Let  $\mathbf{W} \in \mathbb{R}^{2d}$  be the concatenated vector defined as

$$\mathbf{W} = \begin{pmatrix} \mathbf{U} \\ \mathbf{V} \end{pmatrix}. \quad (6)$$

Note that  $\mathbf{W} \xrightarrow{D} N(\mathbf{0}, \Sigma')$  where the covariance matrix is the  $2d$  by  $2d$  block matrix

$$\Sigma' = \begin{bmatrix} \Sigma & 0 \\ 0 & I_d \end{bmatrix} \quad (7)$$

where  $\Sigma$  is given in (3). Since  $\Sigma$  is idempotent, so is  $\Sigma'$ . We next define the  $2d \times 2d$  positive semi-definite matrix  $A$  (composed of four  $d$  by  $d$  block matrices) as

$$A = \begin{bmatrix} I_d & \Lambda \\ \Lambda & \Lambda^2 \end{bmatrix} \quad \text{where} \quad \Lambda = \text{Diag} \left( \frac{\sigma(\epsilon, \delta_n)}{\sqrt{n\mathbf{p}^0}} \right) \quad (8)$$

We can then rewrite our private chi-squared statistic as a quadratic form of the random vectors  $\mathbf{W}$ .

$$Q_{Gauss}^2 = \mathbf{W}^T A \mathbf{W}. \quad (9)$$

*Remark 5.6.* If we have  $\sigma(\epsilon, \delta_n)/\sqrt{n\mathbf{p}^0} \rightarrow \text{constant}$  then the asymptotic distribution of  $Q_{Gauss}^2$  would be a quadratic form of multivariate normals.

Similar to the classical goodness of fit test we consider the limiting case that the random vector  $\mathbf{U}$  is actually a multivariate normal, which will result in  $\mathbf{W}$  being multivariate normal as well. We next want to be able to calculate the distribution of the quadratic form of normals  $\mathbf{W}^T A \mathbf{W}$ . Note that we will write  $\{\chi_1^{2,i}\}_{i=1}^r$  as a set of  $r$  independent chi-squared random variables with one degree of freedom, so that  $\sum_{i=1}^r \chi_1^{2,i} = \chi_r^2$ .

**Theorem 5.7.** *Let  $\mathbf{W} \sim N(\mathbf{0}, \Sigma')$  where  $\Sigma'$  is idempotent and has rank  $r \leq 2d$ . Then the distribution of  $\mathbf{W}^T A \mathbf{W}$  where  $A$  is positive semi-definite is*

$$\sum_{i=1}^r \lambda_i \chi_1^{2,i}$$

where  $\{\lambda_i\}_{i=1}^r$  are the eigenvalues of  $B^T A B$  where  $B \in \mathbb{R}^{2d \times r}$  such that  $BB^T = \Sigma'$  and  $B^T B = I_r$ .

*Proof.* Let  $\mathbf{N}^{(1)} \sim N(\mathbf{0}, I_r)$ . Because  $\Sigma'$  is idempotent, we know that there exists a matrix  $B \in \mathbb{R}^{2d \times r}$  as in the statement of the lemma. Then  $B\mathbf{N}^{(1)}$  has the same distribution as  $\mathbf{W}$ . Also note that because  $B^T A B$  is symmetric, then it is diagonalizable and hence there exists an orthogonal matrix  $H \in \mathbb{R}^{r \times r}$  such that

$$H^T (B^T A B) H = \text{Diag}(\lambda_1, \dots, \lambda_r) \quad \text{where} \quad H^T H = H H^T = I_r$$

Let  $\mathbf{N}^{(1)} = H \mathbf{N}^{(2)}$  where  $\mathbf{N}^{(2)} \sim N(\mathbf{0}, I_r)$ . We then have

$$\mathbf{W}^T A \mathbf{W} \sim (B\mathbf{N}^{(1)})^T A (B\mathbf{N}^{(1)}) \sim (B H \mathbf{N}^{(2)})^T A (B H \mathbf{N}^{(2)}) \sim (\mathbf{N}^{(2)})^T \text{Diag}(\lambda_1, \dots, \lambda_r) \mathbf{N}^{(2)}$$

Now we know that  $(\mathbf{N}^{(2)})^T \text{Diag}(\lambda_1, \dots, \lambda_r) \mathbf{N}^{(2)} \sim \sum_{j=1}^r \lambda_j \chi_1^{2,j}$ , which gives us our result.  $\square$

Note that in the non-private case, the coefficients  $\{\lambda_i\}$  in Theorem 5.7 become the eigenvalues for the rank  $d-1$  idempotent matrix  $\Sigma$ , thus resulting in a  $\chi_{d-1}^2$  distribution. We use the result of Theorem 5.7 in order to find a threshold that will achieve the desired significance level  $1 - \alpha$ , as in the classical chi-squared goodness of fit test. We then set the threshold  $\tau^\alpha$  to satisfy the following:

$$\Pr \left[ \sum_{i=1}^r \lambda_i \chi_1^{2,i} \geq \tau^\alpha \right] = \alpha \tag{10}$$

for  $\{\lambda_i\}$  found in Theorem 5.7. Note, the threshold  $\tau^\alpha$  is a function of  $n, \epsilon, \delta, \alpha$  and  $\mathbf{p}^0$ , but not the data.

We present our modified goodness of fit test when we are dealing with differentially private counts in Algorithm 3.

#### 5.4 Power Analysis of PrivGOF

To determine the power of our new goodness of fit test **PrivGOF**, we need to specify an alternate hypothesis  $H_1 : \mathbf{p} = \mathbf{p}^1$  for  $\mathbf{p}^1 \neq \mathbf{p}^0$ . Similar to past works [Mitra, 1958, Meng and Chapman, 1966, Guenther, 1977], we define parameter  $\tilde{\Delta} > 0$  where

$$\mathbf{p}_n^1 = \mathbf{p}^0 + \frac{\tilde{\Delta}}{\sqrt{n}} (1, -1, \dots, -1, 1)^T \tag{11}$$

---

**Algorithm 3** Private Chi-Squared Goodness of Fit Test

---

```

procedure PRIVGOF(Data  $\mathbf{x}$ ; Privacy  $(\epsilon, \delta)$ , Significance  $1 - \alpha$ ,  $H_0 : \mathbf{p} = \mathbf{p}^0$ )
  Set  $\sigma = \frac{2\sqrt{\log(2/\delta)}}{\epsilon}$ .
  Compute  $Q_{Gauss}^2$  from (4)
  Compute  $\tau^\alpha$  that satisfies (10).
  if  $Q_{Gauss}^2 > \tau^\alpha$  then
    Decision  $\leftarrow$  Reject
  else
    Decision  $\leftarrow$  Fail to Reject
  return Decision.

```

---

for even  $d$ . Note that  $Q_{Gauss}^2$  uses the probability vector given in  $H_0$  but data is generated by Multinomial( $n, \mathbf{p}_n^1$ ). In fact, the nonprivate statistic  $Q^2$  when the data is drawn from  $H_1$  no longer converges to a chi-squared distribution. Instead,  $Q^2$  converges in distribution to a noncentral chi-squared when  $H_1$  holds.<sup>3</sup>

**Lemma 5.8** (Bishop et al. [1975]). *Under the alternate hypothesis  $H_1 : \mathbf{p} = \mathbf{p}_n^1$  given in (11), the chi-squared statistic  $Q^2$  converges in distribution to a noncentral  $\chi_{d-1}^2(\nu)$  where  $\nu = \frac{\tilde{\Delta}^2}{\prod_{i=1}^d p_i^0}$ , i.e. given  $H_1 : \mathbf{p} = \mathbf{p}_n^1$  we have*

$$Q^2 \xrightarrow{D} \chi_{d-1}^2(\nu) \quad \text{as } n \rightarrow \infty.$$

Another classical result tells us that the vector  $\mathbf{U}$  from (2) converges in distribution to a multivariate normal under the alternate hypothesis.

**Lemma 5.9** ([Mitra, 1958, 1955]). *Assume  $\mathbf{X} \sim \text{Multinomial}(n, \mathbf{p}_n^1)$  where  $\mathbf{p}_n^1$  satisfies (11). Then  $\mathbf{U} \xrightarrow{D} N(\mu, \Sigma)$  where  $\Sigma$  is given in (3) and*

$$\mu = \left( \frac{\tilde{\Delta}}{\sqrt{p_1^0}}, \frac{-\tilde{\Delta}}{\sqrt{p_2^0}}, \dots, \frac{-\tilde{\Delta}}{\sqrt{p_{d-1}^0}}, \frac{\tilde{\Delta}}{\sqrt{p_d^0}} \right) \quad (12)$$

**Corollary 5.10.** *Under the alternate hypothesis  $H_1 : \mathbf{p} = \mathbf{p}_n^1$ , then the random vector  $\mathbf{W} \xrightarrow{D} N(\mu', \Sigma')$  for  $\mathbf{W}$  given in (6) where  $\mu' = (\mu, \mathbf{0})^T$  and  $\mu, \Sigma'$  given in (12) and (7), respectively.*

We then write our private statistic as  $Q_{Gauss}^2 = \mathbf{W}^T \mathbf{A} \mathbf{W}$ . Similar to the previous section we will write  $\{\chi_1^{2,i}(\nu_j)\}_{j=1}^r$  as a set of  $r$  independent noncentral chi-squareds with noncentral parameter  $\nu_j$  and one degree of freedom.

**Theorem 5.11.** *Let  $\mathbf{W} \sim N(\mu', \Sigma')$  where  $\mu'$  and  $\Sigma'$  are given in Corollary 5.10. We will write  $\Sigma' = B B^T$  where  $B \in \mathbb{R}^{2d \times (2d-1)}$  has rank  $2d-1$  and  $B^T B = I_{2d-1}$ . We define  $\mathbf{b}^T = (\mu')^T \mathbf{A} B H$*

---

<sup>3</sup>Note that a noncentral chi-squared with noncentral parameter  $\theta$  and  $\nu$  degrees of freedom is the distribution of  $\mathbf{X}^T \mathbf{X}$  where each  $\mathbf{X} \sim N(\mu, I_\nu)$  and  $\theta = \mu^T \mu$ .

where  $H$  is an orthogonal matrix such that  $H^T B^T A B H = \text{Diag}(\lambda_1, \dots, \lambda_{2d-1})$  and  $A$  is given in (8). Then we have

$$\mathbf{W}^T A \mathbf{W} \sim \sum_{j=1}^r \lambda_j \chi_1^{2,j}(\nu_j) + N \left( \kappa, \sum_{j=r+1}^d 4b_j^2 \right), \quad (13)$$

where  $(\lambda_j)_{j=1}^{2d-1}$  are the eigen-values of  $B^T A B$  such that  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0 = \lambda_{r+1} = \dots = \lambda_{2d-1}$  and

$$\nu_j = \left( \frac{b_j}{\lambda_j} \right)^2 \quad \text{for } j \in [r] \quad \& \quad \kappa = \frac{\tilde{\Delta}^2}{\prod_{i=1}^d p_i^0} - \sum_{j=1}^r \frac{b_j^2}{\lambda_j}.$$

*Proof.* We follow a similar analysis as Mohsenipour [2012] for finding the distribution of a quadratic form of normals. Consider the random variable  $\mathbf{N}^{(2)} = B H \mathbf{N}^{(1)} + \mu'$  where  $\mathbf{N}^{(1)} \sim N(\mathbf{0}, I_{2d-1})$ . Note that  $\mathbf{N}^{(2)}$  has the same distribution as  $\mathbf{W}$ . We then have for  $t \geq 0$

$$\begin{aligned} \Pr[\mathbf{W}^T A \mathbf{W} \geq t] &= \Pr[(\mathbf{N}^{(1)})^T H^T B^T A B H \mathbf{N}^{(1)} + 2(\mu')^T A B H \mathbf{N}^{(1)} + (\mu')^T A \mu' \geq t] \\ &= \Pr[(\mathbf{N}^{(1)})^T \text{Diag}(\lambda_1, \dots, \lambda_{d+1}) \mathbf{N}^{(1)} + 2\mathbf{b}^T \mathbf{N}^{(1)} + (\mu')^T A \mu' \geq t] \\ &= \Pr \left[ \sum_{j=1}^r \lambda_j \cdot \left( N_j^{(1)} + b_j / \lambda_j \right)^2 + \sum_{j=r+1}^d 2b_j N_j^{(1)} + \kappa \geq t \right] \\ &= \Pr \left[ \sum_{j=1}^d \lambda_j \cdot \chi_1^{2,j} \left( \left( \frac{b_j}{\lambda_j} \right)^2 \right) + N \left( 0, \sum_{j=r+1}^d 4b_j^2 \right) + \kappa \geq t \right] \end{aligned}$$

□

*Remark 5.12.* Again, if we have  $\sigma(\epsilon, \delta_n) / \sqrt{n \mathbf{p}^0} \rightarrow \text{constant}$  then the asymptotic distribution of  $Q_{Gauss}^2$  converges in distribution to the random variable of the form given in (13) when  $H_1$  from (11) is true.

Obtaining the asymptotic distribution for  $\hat{Q}_{Gauss}^2$  when the alternate hypothesis holds may allow for future results on *effective sample size*, i.e. how large a sample size needs to be in order for **PrivGOF** to have Type II error at most  $\beta$  against  $H_1 : \mathbf{p} = \mathbf{p}_n^1$ . We see this as an important direction for future work.

## 6 Independence Testing

We now consider the problem of testing whether two random variables  $\mathbf{Y}^{(1)} \sim \text{Multinomial}(1, \pi^{(1)})$  and  $\mathbf{Y}^{(2)} \sim \text{Multinomial}(1, \pi^{(2)})$  are independent of each other. Note that  $\sum_{i=1}^r \pi_i^{(1)} = 1$  and  $\sum_{j=1}^c \pi_j^{(2)} = 1$ , so we can write  $\pi_r^{(1)} = 1 - \sum_{i < r} \pi_i^{(1)}$  and  $\pi_c^{(2)} = 1 - \sum_{j < c} \pi_j^{(2)}$ . We then form the null hypothesis  $H_0 : \mathbf{Y}^{(1)} \perp \mathbf{Y}^{(2)}$ , i.e. they are independent. One approach to testing  $H_0$  is to sample  $n$  joint outcomes of  $\mathbf{Y}^{(1)}$  and  $\mathbf{Y}^{(2)}$  and count the number of observed outcomes,  $X_{i,j}$  which is the number of times  $Y_i^{(1)} = 1$  and  $Y_j^{(2)} = 1$  in the  $n$  trials, so that we can summarize all joint outcomes as a contingency table  $\mathbf{X} = (X_{i,j}) \sim \text{Multinomial}(n, \mathbf{p})$ , where  $p_{i,j}$  is the probability that  $Y_i^{(1)} = 1$

and  $Y_j^{(2)} = 1$ . In Table 1 we give a  $r \times c$  contingency table giving the number of joint outcomes for the variables  $\mathbf{Y}^{(1)}$  and  $\mathbf{Y}^{(2)}$  from  $n$  independent trials. We will write the full contingency table of counts  $\mathbf{X} = (X_{i,j})$  as a vector with the ordering convention that we start from the top row and move from left to right across the contingency table.

Table 1: Contingency Table with Marginals.

$Y^{(1)} \backslash Y^{(2)}$	1	2	$\dots$	$c$	Marginals
1	$X_{11}$	$X_{12}$	$\dots$	$X_{1c}$	$X_{1,\cdot}$
2	$X_{21}$	$X_{22}$	$\dots$	$X_{2c}$	$X_{2,\cdot}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$r$	$X_{r,1}$	$X_{r,2}$	$\dots$	$X_{r,c}$	$X_{r,\cdot}$
Marginals	$X_{\cdot,1}$	$X_{\cdot,2}$	$\dots$	$X_{\cdot,c}$	$n$

We want to calculate the chi-squared statistic as in (1) (where now the summation is over all joint outcomes  $i$  and  $j$ ), but now we do not know the true proportion  $\mathbf{p} = (p_{i,j})$  which depends on  $\pi^{(1)}$  and  $\pi^{(2)}$ . However, we can use the *maximum likelihood estimator* (MLE)  $\hat{\mathbf{p}}$  for the probability vector  $\mathbf{p}$  subject to  $H_0$  to form the statistic  $\hat{Q}^2$  where

$$\hat{Q}^2 = \sum_{i,j} \frac{(X_{i,j} - n\hat{p}_{i,j})^2}{n\hat{p}_{i,j}}. \quad (14)$$

The intuition is that if the test rejects even when the most likely probability vector that satisfies the null hypothesis was chosen, then the test should reject against all others. Note that under the null hypothesis we can write  $\mathbf{p}$  as a function of  $\pi^{(1)}$  and  $\pi^{(2)}$ ,

$$\mathbf{p} = \mathbf{f}(\pi^{(1)}, \pi^{(2)}) \text{ where } \mathbf{f} = (f_{i,j}) \quad \text{and} \quad f_{i,j}(\pi^{(1)}, \pi^{(2)}) = \pi_i^{(1)} \cdot \pi_j^{(2)}. \quad (15)$$

Further, we can write the MLE  $\hat{\mathbf{p}}$  as described below.

**Lemma 6.1** ([Bishop et al., 1975]). *Given  $\mathbf{X}$ , which is  $n$  samples of joint outcomes of  $\mathbf{Y}^{(1)} \sim \text{Multinomial}(1, \pi^{(1)})$  and  $\mathbf{Y}^{(2)} \sim \text{Multinomial}(1, \pi^{(2)})$ , if  $\mathbf{Y}^{(1)} \perp \mathbf{Y}^{(2)}$ , then the MLE for  $\mathbf{p} = \mathbf{f}(\pi^{(1)}, \pi^{(2)})$  for  $\mathbf{f}$  given in (15) is the following:  $\hat{\mathbf{p}} = \mathbf{f}(\hat{\pi}^{(1)}, \hat{\pi}^{(2)})$  where*

$$\hat{\pi}_i^{(1)} = X_{i,\cdot}/n, \quad \hat{\pi}_j^{(2)} = X_{\cdot,j}/n \quad \text{for } i \in [r], j \in [c] \quad (16)$$

where  $X_{i,\cdot} = \sum_{j=1}^c X_{i,j}$  and  $X_{\cdot,j} = \sum_{i=1}^r X_{i,j}$ .

We then state another classical result that gives the asymptotic distribution of  $\hat{Q}^2$  given  $H_0$ .

**Theorem 6.2.** *Bishop et al. [1975] Given the assumptions in Lemma 6.1, the statistic  $\hat{Q}^2$  given in (14) converges in distribution to a chi-squared distribution, i.e.*

$$\hat{Q}^2 \xrightarrow{D} \chi_\nu^2$$

for  $\nu = (r-1)(c-1)$ .



---

**Algorithm 4** Pearson Chi-Squared Independence Test

---

```

procedure INDEP(Data  $\mathbf{x}$ , and Significance  $1 - \alpha$ )
   $\hat{\mathbf{p}} \leftarrow$  MLE calculation in (16)
  Compute  $\hat{Q}^2$  from (14) and set  $\nu = (r - 1)(c - 1)$ .
  if  $\hat{Q}^2 > \chi_{\nu, 1-\alpha}^2$  and all entries of  $\mathbf{x}$  are at least 5 then
    Decision  $\leftarrow$  Reject
  else
    Decision  $\leftarrow$  Fail to Reject
  return Decision.

```

---

The chi-squared independence test is then to compare the statistic  $\hat{Q}^2$ , with the value  $\chi_{\nu, 1-\alpha}^2$  for a  $1 - \alpha$  significance test. We formally give the Pearson Chi-Squared test in Algorithm 4. An often used “rule of thumb” [Triola, 2014] with this test is that it can only be used if all the cell counts are at least 5, otherwise the test Fails to Reject  $H_0$ . We will follow this rule of thumb in our tests.

Similar to our prior analysis for goodness of fit, we aim to understand the asymptotic distribution from Theorem 6.2. First, we can define  $\hat{\mathbf{U}}$  in terms of the MLE  $\hat{\mathbf{p}}$  given in (16):

$$\hat{U}_{i,j} = (X_{i,j} - n\hat{p}_{i,j}) / \sqrt{n\hat{p}_{i,j}}. \quad (17)$$

The following classical result gives the asymptotic distribution of  $\hat{\mathbf{U}}$  under  $H_0$ , which also proves Theorem 6.2.

**Lemma 6.3.** [Bishop et al., 1975] *With the same hypotheses as Lemma 6.1, the random vector  $\hat{\mathbf{U}}$  given in (17) converges in distribution to a multivariate normal,*

$$\hat{\mathbf{U}} \xrightarrow{D} N(0, \Sigma_{ind})$$

where  $\Sigma_{ind} = I_{rc} - \sqrt{\mathbf{p}} \cdot \sqrt{\mathbf{p}}^T - \Gamma(\Gamma^T \Gamma)^{-1} \Gamma^T$  with  $\mathbf{f}$  given in (15), and

$$\Gamma = \text{Diag}(\sqrt{\mathbf{p}})^{-1} \cdot \nabla \mathbf{f}(\pi^{(1)}, \pi^{(2)}),$$

$$\nabla \mathbf{f}(\pi^{(1)}, \pi^{(2)}) = \begin{bmatrix} \frac{\partial f_{1,1}}{\partial \pi_1^{(1)}} & \cdots & \frac{\partial f_{1,1}}{\partial \pi_{r-1}^{(1)}} & \frac{\partial f_{1,1}}{\partial \pi_1^{(2)}} & \cdots & \frac{\partial f_{1,1}}{\partial \pi_{c-1}^{(2)}} \\ \frac{\partial f_{1,2}}{\partial \pi_1^{(1)}} & \cdots & \frac{\partial f_{1,2}}{\partial \pi_{r-1}^{(1)}} & \frac{\partial f_{1,2}}{\partial \pi_1^{(2)}} & \cdots & \frac{\partial f_{1,2}}{\partial \pi_{c-1}^{(2)}} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_{r,c}}{\partial \pi_1^{(1)}} & \cdots & \frac{\partial f_{r,c}}{\partial \pi_{r-1}^{(1)}} & \frac{\partial f_{r,c}}{\partial \pi_1^{(2)}} & \cdots & \frac{\partial f_{r,c}}{\partial \pi_{c-1}^{(2)}} \end{bmatrix}_{rc, r+c-2}.$$

In order to do a test that is similar to **Indep** given in Algorithm 4, we need to determine an estimate for  $\pi^{(1)}$  and  $\pi^{(2)}$  where we are only given access to the noisy cell counts.

## 6.1 Estimating Parameters with Private Counts

We now assume that we do not have access to the counts  $X_{i,j}$  from Table 1 but instead we have  $W_{i,j} = X_{i,j} + Z_{i,j}$  where  $Z_{i,j} \sim \mathcal{D}$  for Laplace or Gaussian noise given in (5) and we want to perform a test for independence. We consider the full likelihood of the noisy  $r \times c$  contingency table

$$\begin{aligned} \Pr [\mathbf{X} + \mathbf{Z} = \mathbf{w} | H_0, \pi^{(1)}, \pi^{(2)}] &= \sum_{\substack{\mathbf{x}: \sum_{i,j} x_{i,j} = n \\ x_{i,j} \in \mathbb{N}}} \Pr [\mathbf{X} = \mathbf{x} | \pi^{(1)}, \pi^{(2)}] \cdot \Pr [\mathbf{Z} = \mathbf{w} - \mathbf{x} | \mathbf{X} = \mathbf{x}] \\ &= \sum_{\substack{\mathbf{x}: \sum_{i,j} x_{i,j} = n \\ x_{i,j} \in \mathbb{N}}} \underbrace{\Pr [\mathbf{X} = \mathbf{x} | \pi^{(1)}, \pi^{(2)}]}_{\text{Multinomial}} \prod_{i,j} \underbrace{\Pr [Z_{i,j} = w_{i,j} - X_{i,j} | \mathbf{X} = \mathbf{x}]}_{\text{Noise}} \end{aligned}$$

to find the best estimates for  $\{\pi^{(i)}\}$  given the noisy counts.

---

### Algorithm 5 Two Step MLE Calculation

---

**procedure** 2MLE(Noisy Data  $\mathbf{X} + \mathbf{Z} = \mathbf{w}$ )

$\tilde{\mathbf{x}} \leftarrow$  Solution to (18). If  $\mathcal{D} = \text{Gauss}$ , then  $\gamma = 1$ , if  $\mathcal{D} = \text{Lap}$  set  $0 < \gamma \ll 1$ .

**if** Any cell of  $\tilde{\mathbf{x}}$  is less than 5 **then**

$\tilde{\pi}^{(1)}, \tilde{\pi}^{(2)} \leftarrow \text{NULL}$

**else**

$\tilde{\pi}^{(1)}, \tilde{\pi}^{(2)} \leftarrow$  MLE for  $\pi^{(1)}$  and  $\pi^{(2)}$  with data  $\tilde{\mathbf{x}}$ .

Note that the MLE for the probabilities  $\pi^{(1)}$  and  $\pi^{(2)}$  with data is given in (16).

**return**  $\tilde{\pi}^{(1)}$  and  $\tilde{\pi}^{(2)}$ .

---

Maximizing this quantity is computationally very expensive for values of  $n > 100$  even for  $2 \times 2$  tables,<sup>4</sup> so we instead follow a two step procedure similar to the work of Karwa and Slavković [2016], where they “denoise” a private degree sequence for a synthetic graph and then use the denoised estimator to approximate the parameters of the  $\beta$ -model of random graphs. We will first find the most likely contingency table given the noisy data  $\mathbf{w}$  and then find the most likely probability vectors under the null hypothesis that could have generated that denoised contingency table (this is not equivalent to maximizing the full likelihood, but it seems to work well as our experiments later show). For the latter step, we use Equation (16) to get the MLE for  $\pi^{(1)}$  and  $\pi^{(2)}$  given a vector of counts  $\mathbf{x}$ . For the first step, we need to minimize  $\|\mathbf{w} - \mathbf{x}\|$  subject to  $\sum_{i,j} x_{i,j} = n$  and  $x_{i,j} \geq 0$  where the norm in the objective is either  $\ell_1$  for Laplace noise or  $\ell_2$  for Gaussian noise.

Note that for Laplace noise, the above optimization problem does not give a unique solution and it is not clear which contingency table  $\mathbf{x}$  to use. One solution to overcome this is to add a *regularizer* to the objective value. We will follow the work of Lee et al. [2015] to overcome this problem by using an *elastic net* regularizer [Zou and Hastie, 2005]:

$$\begin{aligned} \underset{\mathbf{x}}{\text{argmin}} \quad & (1 - \gamma) \cdot \|\mathbf{w} - \mathbf{x}\|_1 + \gamma \cdot \|\mathbf{w} - \mathbf{x}\|_2^2 \\ \text{s.t.} \quad & \sum_{i,j} x_{i,j} = n, \quad x_{i,j} \geq 0. \end{aligned} \tag{18}$$

---

<sup>4</sup>Note that there is a  $\text{poly}(n)$  time algorithm to solve this, but the coefficients in each term of the sum can be very large numbers, with  $\text{poly}(n)$  bits, which makes it difficult for numeric solvers.

where if we use Gaussian noise, we set  $\gamma = 1$  and if we use Laplace noise then we pick a small  $\gamma > 0$  and then solve the resulting program. Our two step procedure for finding an approximate MLE for  $\pi^{(1)}$  and  $\pi^{(2)}$  based on our noisy vector of counts  $\mathbf{w}$  is given in Algorithm 5, where we take into account the rule of thumb from **Indep** and return NULL if any computed table has counts less than 5.

We will denote  $\tilde{\mathbf{p}}$  to be the probability vector of  $\mathbf{f}$  from (15) applied to the result of  $2\text{MLE}(\mathbf{X} + \mathbf{Z})$ . We now write down the private chi-squared statistic when we use the estimate  $\tilde{\mathbf{p}}$  in place of the actual (unknown) probability vector  $\mathbf{p}$ :

$$\tilde{Q}_{\mathcal{D}}^2 = \sum_{i,j} \frac{(X_{i,j} + Z_{i,j} - n\tilde{p}_{i,j})^2}{n\tilde{p}_{i,j}} \quad \{Z_{i,j}\} \stackrel{i.i.d.}{\sim} \mathcal{D}. \quad (19)$$

---

**Algorithm 6** MC Independence Testing

---

**procedure**  $\text{MCIndep}_{\mathcal{D}}$ (Contingency Table  $\mathbf{x}$ ; privacy parameters  $(\epsilon, \delta)$ , significance  $1 - \alpha$ )

$\mathbf{w} \leftarrow \mathbf{x} + \mathbf{Z}$ , where  $\{Z_{i,j}\} \stackrel{i.i.d.}{\sim} \mathcal{D}$  and  $\mathcal{D}$  given in (5).

$(\tilde{\pi}^{(1)}, \tilde{\pi}^{(2)}) \leftarrow 2\text{MLE}(\mathbf{w})$ .

**if**  $(\tilde{\pi}^{(1)}, \tilde{\pi}^{(2)}) == \text{NULL}$  **then return** Fail to Reject

**else**

$\tilde{q} \leftarrow \tilde{Q}_{\mathcal{D}}^2$  with data  $\mathbf{w}$  and parameters  $\tilde{\mathbf{p}} = \mathbf{f}(\tilde{\pi}^{(1)}, \tilde{\pi}^{(2)})$ .

Set  $k > 1/\alpha$  and  $q \leftarrow \text{NULL}$ .

**for**  $t \in [k]$  **do**

Generate a fresh contingency table  $\tilde{\mathbf{x}}$  using parameters  $(\tilde{\pi}^{(1)}, \tilde{\pi}^{(2)})$ .

$\tilde{\mathbf{w}} \leftarrow \tilde{\mathbf{x}} + \mathbf{Z}$ , where  $\{Z_{i,j}\} \stackrel{i.i.d.}{\sim} \mathcal{D}$  and  $\mathcal{D}$  given in (5).

$\tilde{\pi}^1, \tilde{\pi}^2 \leftarrow 2\text{MLE}(\tilde{\mathbf{w}})$ .

**if**  $\tilde{\pi}^1, \tilde{\pi}^2 == \text{NULL}$  **then return** Fail to Reject

**else**

Concatenate  $q$  with  $\tilde{Q}_{\mathcal{D}}^2$  given in (19) with  $\tilde{\mathbf{w}}$  and  $\tilde{\mathbf{p}} = \mathbf{f}(\tilde{\pi}^1, \tilde{\pi}^2)$ .

$\tilde{\tau}^\alpha \leftarrow$  the  $\lceil (k+1)(1-\alpha) \rceil$  ranked statistic in  $q$ .

**if**  $\tilde{q} > \tilde{\tau}^2$  **then return** Reject  $H_0$ .

**else return** Fail to Reject  $H_0$ .

---

## 6.2 Monte Carlo Test: $\text{MCIndep}_{\mathcal{D}}$

We first follow a similar procedure as in Section 5.2 but using the parameter estimates from  $2\text{MLE}$  instead of the actual (unknown) probabilities. Our procedure  $\text{MCIndep}_{\mathcal{D}}$  (given in Algorithm 6) works as follows: given a dataset  $\mathbf{x}$ , we will add the appropriately scaled Laplace or Gaussian noise to ensure differential privacy to get the noisy table  $\mathbf{w}$ . Then we use  $2\text{MLE}$  on the private data to get approximates to the parameters  $\pi^{(i)}$ , which we denote as  $\tilde{\pi}^{(i)}$  for  $i = 1, 2$ . Using these probability estimates, we sample  $k > 1/\alpha$  many contingency tables and noise terms to get  $k$  different values for  $\tilde{Q}_{\mathcal{D}}^2$  and choose the  $\lceil (k+1)(1-\alpha) \rceil$  ranked statistic as our threshold  $\tilde{\tau}^\alpha$ . If at any stage  $2\text{MLE}$  returns NULL, then the test Fails to Reject  $H_0$ . We formally give our test  $\text{MCIndep}_{\mathcal{D}}$  in Algorithm 6.

### 6.3 Asymptotic Approach: PrivIndep

---

**Algorithm 7** Private Independence Test for  $r \times c$  tables

---

```

procedure PRIVINDEP(Data  $\mathbf{x}$ , privacy parameters  $(\epsilon, \delta)$ , and Significance  $1 - \alpha$ )
    Compute the private contingency table  $\mathbf{X} + \mathbf{Z} = \mathbf{w}$  where  $\mathbf{Z} \sim N(0, \sigma^2 I_{r \cdot c})$  and  $\sigma$  from (5).
     $(\tilde{\pi}^{(1)}, \tilde{\pi}^{(2)}) \leftarrow \text{2MLE}(\mathbf{w})$ .
    if  $(\tilde{\pi}^{(1)}, \tilde{\pi}^{(2)}) == \text{NULL}$  then
        Decision  $\leftarrow$  Fail to Reject
    else
         $\tilde{\mathbf{p}} \leftarrow \mathbf{f}(\tilde{\pi}^{(1)}, \tilde{\pi}^{(2)})$  for  $\mathbf{f}$  given in (15).
        Compute  $\tilde{Q}_{Gauss}^2$  from (19) with noisy data  $\mathbf{w}$ .
        Compute  $\tilde{\tau}^\alpha$  that satisfies (22).
        if  $\tilde{Q}_{Gauss}^2 > \tilde{\tau}^\alpha$  then
            Decision  $\leftarrow$  Reject
        else
            Decision  $\leftarrow$  Fail to Reject
    return Decision.

```

---

We will now focus on the analytical form of our private statistic when Guassian noise is added. We can then write  $\tilde{Q}_{Gauss}^2$  in its quadratic form, which is similar to the form of  $Q_{Gauss}^2$  from (9),

$$\tilde{Q}_{Gauss}^2 = \tilde{\mathbf{W}}^T \tilde{A} \tilde{\mathbf{W}} \quad (20)$$

where  $\tilde{\mathbf{W}} = (\tilde{\mathbf{U}})$  with  $\tilde{\mathbf{U}}$  set in (17) except with  $\tilde{\mathbf{p}}$  used instead of the given  $\mathbf{p}^0$  in the goodness of fit testing and  $\mathbf{V}$  set as in (6). Further, we denote  $\tilde{A}$  as  $A$  in (8) but with  $\tilde{\mathbf{p}}$  instead of  $\mathbf{p}^0$ . We will use the  $2rc$  by  $2rc$  block matrix  $\tilde{\Sigma}'_{ind}$  to estimate the covariance of  $\tilde{\mathbf{W}}$ , where

$$\tilde{\Sigma}'_{ind} = \begin{bmatrix} \tilde{\Sigma}_{ind} & 0 \\ 0 & I_{rc} \end{bmatrix} \quad (21)$$

and  $\tilde{\Sigma}_{ind}$  is the matrix  $\Sigma_{ind}$  in Lemma 6.3, except we use our estimates  $\tilde{\pi}^{(1)}$ ,  $\tilde{\pi}^{(2)}$ , or  $\tilde{\mathbf{p}}$  whenever we need to use the actual (unknown) parameters.

Thus, if we are given a differentially private version of a contingency table where each cell has added independent Gaussian noise with variance  $\sigma^2$ , we calculate  $\tilde{Q}_{Gauss}^2$  and compare it to the threshold  $\tilde{\tau}^\alpha$  where

$$\Pr \left[ \sum_{i=1}^{rc} \tilde{\lambda}_i \chi_1^{2,i} \geq \tilde{\tau}^\alpha \right] = \alpha \quad (22)$$

with  $\{\tilde{\lambda}_i\}$  being the eigenvalues of  $\tilde{B}^T \tilde{A} \tilde{B}$  with rank  $\nu = rc + (r-1)(c-1)$  matrix  $\tilde{B} \in \mathbb{R}^{2rc, \nu}$  where  $\tilde{B} \tilde{B}^T = \tilde{\Sigma}'_{ind}$ . Our new independence test **PrivIndep** is given in Algorithm 7, where **2MLE** estimates  $\pi^{(i)}$  for  $i = 1, 2$  and **PrivIndep** Fails to Reject if **2MLE** returns NULL.

## 7 Significance Results

We now show how each of our tests perform on simulated data when  $H_0$  holds in goodness of fit and independence testing. We fix our desired significance  $1 - \alpha = 0.95$  and privacy level

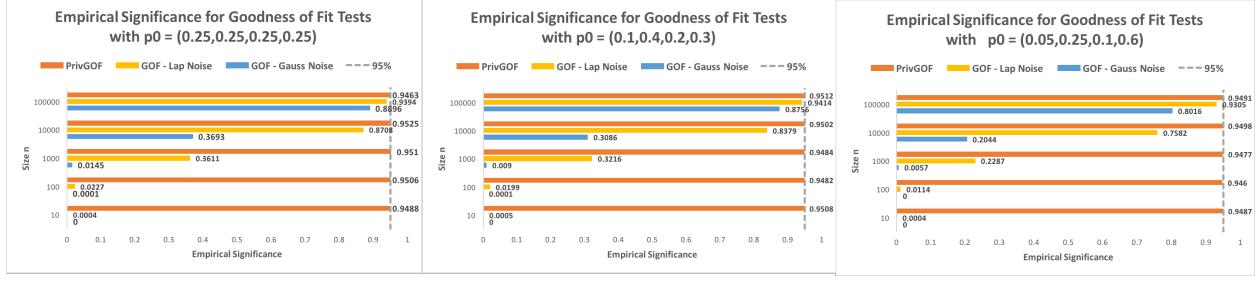


Figure 1: Significance of the classical test **GOF** when used on counts with added Laplace or Gaussian noise compared to **PrivGOF** in 10,000 trials with  $(\epsilon, \delta) = (0.1, 10^{-6})$  and  $\alpha = 0.05$ .

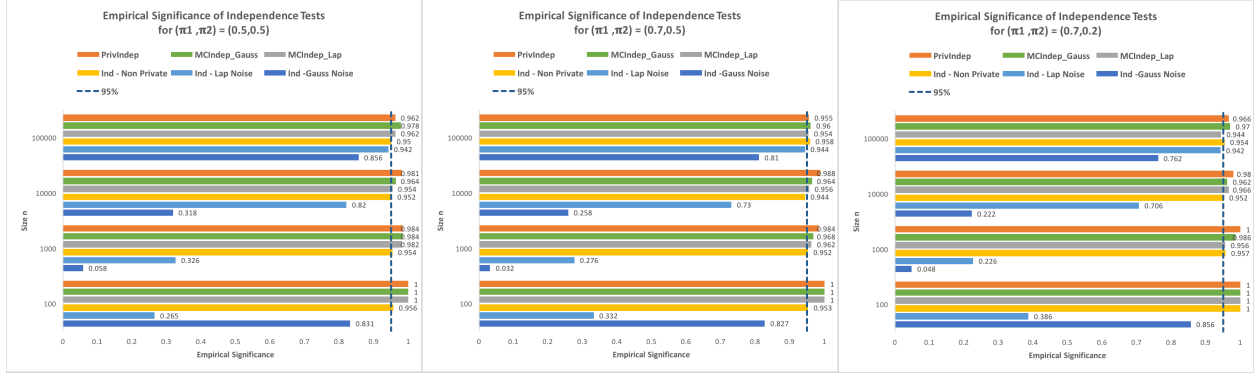


Figure 2: Significance of **Indep** when used on a contingency table with added Laplace or Gaussian noise compared to **MCIndep<sub>D</sub>** for both Laplace and Gaussian noise and **PrivIndep** in 1,000 trials with  $(\epsilon, \delta) = (0.1, 10^{-6})$  and  $\alpha = 0.05$ .

$(\epsilon, \delta) = (0.1, 10^{-6})$  in all of our tests.

By Theorem 5.5, we know that **MCGOF<sub>D</sub>** will have significance at least  $1 - \alpha$ . We then turn to our test **PrivGOF** to compute the proportion of trials that failed to reject  $H_0 : \mathbf{p} = \mathbf{p}^0$  when it holds. In Figure 1 we give several different null hypotheses  $\mathbf{p}^0$  and sample sizes  $n$  to show that **PrivGOF** achieves near 0.95 significance in all our tested cases. We also compare our results with how the original test **GOF** would perform if used on the private counts with either Laplace and Gaussian noise.

To show that **PrivGOF** works beyond  $d = 4$  Multinomial data, we give a table of results in Table 2 for  $d = 100$  data and null hypothesis  $p_i^0 = 1/100$  for  $i \in [100]$ . We give the proportion of 10,000 trials that were not rejected by **PrivGOF** in the “**PrivGOF Sign**” column and those that were not rejected by the classical test **GOF** in the “**Indep Sign**” column. Note that the critical value that **GOF** uses is 123.23 for every test in this case, whereas **PrivGOF**’s critical value changes for each test.

We then turn to independence testing for  $2 \times 2$  contingency tables using both **MCIndep<sub>D</sub>** and **PrivIndep**. Note that our methods do apply to arbitrary  $k \times \ell$  tables and run in time  $\text{poly}(k, \ell, \log(n))$  plus the time for the iterative Imhof method to find the critical values. In Figure 2 we compute the empirical significance of both of our tests and compare it to how **Indep** performs

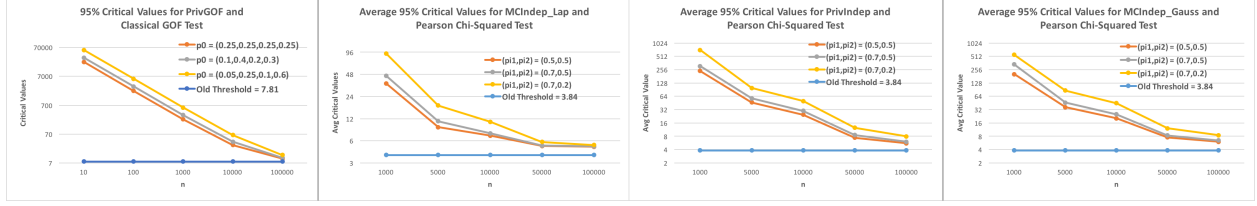


Figure 3: Comparison of the (average) critical values for all of our tests with  $\alpha = 0.05$  and  $(\epsilon, \delta) = (0.1, 10^{-6})$ . Note that some are on a log-scale.

Table 2: Goodness of fit testing for multinomial data with  $\alpha = 0.05$  and  $(\epsilon, \delta) = (0.1, 10^{-6})$  for dimension  $d = 100$  data.

$\mathbf{p}^0$			$n$	$\chi^2_{d-1, 1-\alpha}$	Indep Signf	$\tau^\alpha$	PrivGOF Signf
0.01	...	0.01	1,500	123.23	0.0000	48,231	0.9522
0.01	...	0.01	10,000	123.23	0.0000	7,339	0.9491
0.01	...	0.01	100,000	123.23	0.0000	844.7	0.9511
0.01	...	0.01	1,000,000	123.23	0.0524	195.3	0.9479

on the nonprivate data. For **MCIndep<sub>D</sub>** and **PrivIndep** we sample 1,000 trials for various parameters  $\pi^{(1)}$ ,  $\pi^{(2)}$ , and  $n$  that could have generated the contingency tables. We set the number of samples  $k = 50$  in **MCIndep<sub>D</sub>** regardless of the noise we added and when we use Laplace noise, we set  $\gamma = 0.01$  as the parameter in **2MLE**. Note that when  $n$  is small, we get that our differentially private independence tests almost always fails to reject. In fact, when  $n = 100$  all of our tests in 1,000 trials fail to reject. This is due to **2MLE** releasing a contingency table based on the private counts with small cell counts. When the cell counts in **2MLE** are small we follow the “rule of thumb” from the classical test **Indep** and output NULL, which results in **PrivIndep** failing to reject. This will ensure good significance but makes no promises on power for small  $n$ , as does the classical test **Indep**. Further, another consequence of this “rule of thumb” is that when we use **Indep** on private counts, with either Laplace or Gaussian noise, it tends to have lower Type I error than for larger  $n$ .

We also plot the critical values of our various tests in Figure 3. For both **PrivGOF** and **PrivIndep** we used the package in R “**CompQuadForm**” that has various methods for finding estimates to the tail probabilities for quadratic forms of normals, of which we used the “**imhof**” method [Imhof, 1961] to approximate the threshold for each test. Note that in **MCIndep<sub>D</sub>** and **PrivIndep** each trial has a different threshold, so we give the average over all trials.

## 8 Power Results

We now want to show that our tests correctly reject  $H_0$  when it is false, fixing parameters  $\alpha = 0.05$  and  $(\epsilon, \delta) = (0.1, 10^{-6})$ . For our two goodness of fit tests, **MCGOF<sub>D</sub>** (with  $k = 100$ ) and **PrivGOF** we test whether the multinomial data came from  $\mathbf{p}^0 = (1/4, 1/4, 1/4, 1/4)$  when it was actually sampled from  $\mathbf{p}^1 = \mathbf{p}^0 + 0.01 \cdot (1, -1, 1, -1)$ . We compare each of our tests with the classical **Indep** test that uses the unaltered data in Figure 4. We then find the proportion of 1,000 trials that each

of our tests rejected  $H_0 : \mathbf{p} = \mathbf{p}^0$  for various  $n$ . Note that **Indep** has difficulty distinguishing  $\mathbf{p}^0$  and  $\mathbf{p}^1$  for reasonable sample sizes.

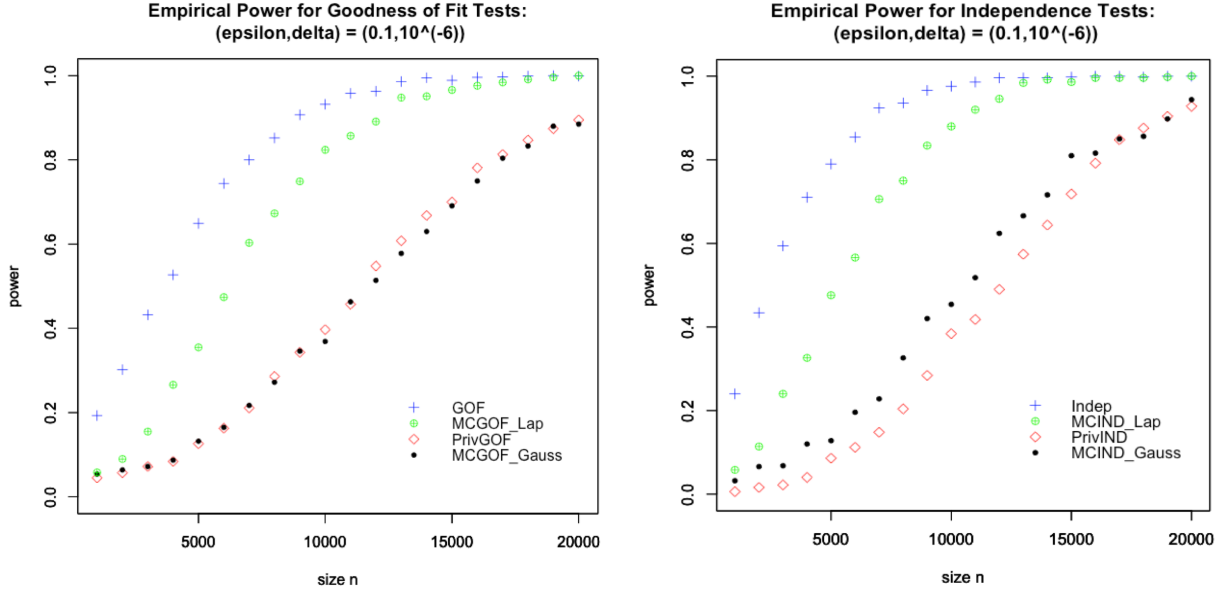


Figure 4: Power Plots of  $\text{MCGOF}_{\mathcal{D}}$  and  $\text{PrivGOF}$  with alternate  $\mathbf{p}^1$  with parameter  $\Delta = 0.01$ , as well as our independence tests  $\text{MCIndep}_{\mathcal{D}}$  and  $\text{PrivIndep}$  compared with the classical tests with alternate covariance 0.01, with  $(\epsilon, \delta) = (0.1, 10^{-6})$ .

We then turn to independence testing for  $2 \times 2$  tables with our two differentially private tests  $\text{MCIndep}_{\mathcal{D}}$  and  $\text{PrivIndep}$ . We fix the alternate  $H_1 : \text{Cov}(\mathbf{Y}^{(1)}, \mathbf{Y}^{(2)}) = \Delta > 0$  so that  $\mathbf{Y}^{(1)} \sim \text{Bern}(\pi^{(1)} = 1/2)$  and  $\mathbf{Y}^{(2)} \sim \text{Bern}(\pi^{(2)} = 1/2)$  are not independent. We then sample contingency tables from a multinomial distribution with probability  $\mathbf{p}^1 = (1/4, 1/4, 1/4, 1/4) + \Delta(1, -1, 1, -1)$  and various sizes  $n$ . We compute the proportion of 1,000 trials that  $\text{MCIndep}_{\mathcal{D}}$  and  $\text{PrivIndep}$  rejected  $H_0 : \mathbf{Y}^{(1)} \perp \mathbf{Y}^{(2)}$  and  $\Delta = 0.01$  in Figure 4. For  $\text{MCIndep}_{\mathcal{D}}$  we set the number of samples  $k = 50$  and when we use Laplace noise, we set  $\gamma = 0.01$  in 2MLE.

## 9 Conclusion

We proposed new hypothesis tests based on a private version of the chi-squared statistic for goodness of fit and independence tests. For each test, we showed analytically or experimentally that we can achieve significance close to the target  $1 - \alpha$  level similar to the nonprivate tests. We also showed that all the tests have a loss in power with respect to the non-private classical tests, with methods using Laplace noise outperforming those with Gaussian noise, due to the fact that the Gaussian noise has higher variance (to achieve the same level of privacy). Experimentally we show for  $2 \times 2$  tables that with less than 3000 additional samples the tests with Laplace noise achieve the same power as the classical tests. Typically, one would expect differential privacy to require the sample size to blow up by a multiplicative  $1/\epsilon$  factor. However, we see a better performance because the noise is dominated by the sampling error.

## Acknowledgements

We would like to thank the following people for helpful discussions: Dan Kifer, Aaron Roth, Aleksandra B. Slavković, Or Sheffet, Adam Smith, and a number of others involved with the Privacy Tools for Sharing Research Data project. A special thanks to Vishesh Karwa for helping us with statistics background and the suggestion to use a two step MLE procedure.



## References

- Boaz Barak, Kamalika Chaudhuri, Cynthia Dwork, Satyen Kale, Frank McSherry, and Kunal Talwar. Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In *Proceedings of the Twenty-Sixth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, June 11-13, 2007, Beijing, China*, pages 273–282, 2007.
- Brian JL Berry. City size distributions and economic development. *Economic development and cultural change*, pages 573–588, 1961.
- Yvonne M. M. Bishop, Stephen E. Fienberg, and Paul W. Holland. Discrete multivariate analysis: Theory and practice, 1975.
- Aaron Blair, Pierre Decoufle, and D Grauman. Causes of death among laundry and dry cleaning workers. *American journal of public health*, 69(5):508–511, 1979.
- TJ David and SC Beards. Asthma and the month of birth. *Clinical & Experimental Allergy*, 15(4):391–395, 1985.
- Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our data, ourselves: Privacy via distributed noise generation. In *Proceedings of the 24th Annual International Conference on The Theory and Applications of Cryptographic Techniques*, EURO-CRYPT’06, 2006a.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *TCC ’06*, pages 265–284, 2006b.
- Helen Rose Fuchs Ebaugh and C Allen Haney. Church attendance and attitudes toward abortion: Differentials in liberal and conservative churches. *Journal for the Scientific Study of Religion*, pages 407–413, 1978.
- Stephen E. Fienberg, Alessandro Rinaldo, and Xiaolin Yang. Differential privacy and the risk-utility tradeoff for multi-dimensional contingency tables. In *Proceedings of the 2010 International Conference on Privacy in Statistical Databases*, PSD’10, pages 187–199, Berlin, Heidelberg, 2010. Springer-Verlag.
- Marco Gaboardi, Emilio Jesús Gallego Arias, Justin Hsu, Aaron Roth, and Zhiwei Steven Wu. Dual query: Practical private query release for high dimensional data. In *International Conference on Machine Learning*, 2014.
- J Cox Gill, Janet Endres-Brooks, Patricia J Bauer, William J Marks Jr, and Robert R Montgomery. The effect of abo blood group on the diagnosis of von willebrand disease. *Blood*, 69(6):1691–1695, 1987.
- William A Glaser. The family and voting turnout. *Public Opinion Quarterly*, 23(4):563–570, 1959.
- Anthony G Greenwald, Catherine G Carnot, Rebecca Beach, and Barbara Young. Increasing voting behavior by asking people if they expect to vote. *Journal of Applied Psychology*, 72(2):315, 1987.
- William C. Guenther. Power and sample size for approximate chi-square tests. *The American Statistician*, 31(2):83–85, 1977.

- Moritz Hardt, Katrina Ligett, and Frank McSherry. A simple and practical algorithm for differentially private data release. In *Conference on Neural Information Processing Systems (NIPS)*, pages 2348–2356, 2012.
- Nils Homer, Szabolcs Szelinger, Margot Redman, David Duggan, Waibhav Tembe, Jill Muehling, John V. Pearson, Dietrich A. Stephan, Stanley F. Nelson, and David W. Craig. Resolving individuals contributing trace amounts of dna to highly complex mixtures using high-density snp genotyping microarrays. *PLoS Genet*, 4(8), 08 2008.
- J. P. Imhof. Computing the distribution of quadratic forms in normal variables. *Biometrika*, 48 (3-4):419–426, 1961.
- Aaron Johnson and Vitaly Shmatikov. Privacy-preserving data exploration in genome-wide association studies. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’13, pages 1079–1087, New York, NY, USA, 2013. ACM.
- Vishesh Karwa and Aleksandra Slavković. Inference using noisy degrees: Differentially private  $\beta$ -model and synthetic graphs. *Ann. Statist.*, 44(P1):87–112, 02 2016.
- Vishesh Karwa and Aleksandra Slavković. Differentially private graphical degree sequences and synthetic graphs. In Josep Domingo-Ferrer and Ilenia Tinnirello, editors, *Privacy in Statistical Databases*, volume 7556 of *Lecture Notes in Computer Science*, pages 273–285. Springer Berlin Heidelberg, 2012.
- Matthew Krain and Marissa Edson Myers. Democracy and civil war: A note on the democratic peace proposition. *International Interactions*, 23(1):109–118, 1997.
- James H Kuklinski and Darrell M West. Economic expectations and voting behavior in united states house and senate elections. *American Political Science Review*, 75(02):436–447, 1981.
- Jaewoo Lee, Yue Wang, and Daniel Kifer. Maximum likelihood postprocessing for differential privacy under consistency constraints. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’15, pages 635–644, New York, NY, USA, 2015. ACM.
- Chao Li and Gerome Miklau. An adaptive mechanism for accurate query answering under differential privacy. *Proc. VLDB Endow.*, 5(6):514–525, February 2012.
- Chao Li, Michael Hay, Vibhor Rastogi, Gerome Miklau, and Andrew McGregor. Optimizing linear counting queries under differential privacy. In *ACM SIGACT–SIGMOD–SIGART Symposium on Principles of Database Systems (PODS)*, pages 123–134, 2010.
- Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science*, FOCS ’07, pages 94–103, Washington, DC, USA, 2007. IEEE Computer Society.
- Rosa C. Meng and Douglas G. Chapman. The power of chi square tests for contingency tables. *Journal of the American Statistical Association*, 61(316):965–975, 1966.
- Neil J Mitchell and James M McCormick. Economic and political explanations of human rights violations. *World Politics*, 40(04):476–498, 1988.

- S.K. Mitra. *Contributions to the Statistical Analysis of Categorical Data*. Institute of Statistics mimeo series. 1955.
- S.K. Mitra. On the limiting power function of the frequency chi-square test. *Ann. Math. Statist.*, 29(4):1221–1233, 12 1958.
- Ali Akbar Mohsenipour. *On the Distribution of Quadratic Expressions in Various Types of Random Vectors*. PhD thesis, The University of Western Ontario, Electronic Thesis and Dissertation Repository, 12 2012. Paper 955.
- Sean Simmons and Bonnie Berger. Realizing privacy preserving genome-wide association studies. *Bioinformatics*, 32(9):1293–1300, 2016. doi: 10.1093/bioinformatics/btw009. URL <http://bioinformatics.oxfordjournals.org/content/32/9/1293.abstract>.
- Adam Smith. Privacy-preserving statistical estimation with optimal convergence rates. In *Proceedings of the Forty-third Annual ACM Symposium on Theory of Computing*, STOC ’11, pages 813–822, New York, NY, USA, 2011. ACM.
- M.F. Triola. *Essentials of Statistics*. Pearson Education, 2014. ISBN 9780321924636. URL <https://books.google.com/books?id=QZN-AgAAQBAJ>.
- Caroline Uhler, Aleksandra Slavkovic, and Stephen E. Fienberg. Privacy-preserving data sharing for genome-wide association studies. *Journal of Privacy and Confidentiality*, 5(1), 2013.
- Duy Vu and Aleksandra Slavkovic. Differential privacy for clinical trial data: Preliminary evaluations. In *Proceedings of the 2009 IEEE International Conference on Data Mining Workshops*, ICDMW ’09, pages 138–143, Washington, DC, USA, 2009. IEEE Computer Society.
- Y. Wang, J. Lee, and D. Kifer. Differentially Private Hypothesis Testing, Revisited. *ArXiv e-prints*, November 2015.
- Larry Wasserman and Shuheng Zhou. A statistical framework for differential privacy. *Journal of the American Statistical Association*, 105(489):375–389, 2010.
- Fei Yu and Zhanglong Ji. Scalable privacy-preserving data sharing methodology for genome-wide association studies: an application to idash healthcare privacy protection challenge. *BMC Medical Informatics and Decision Making*, 14(1):1–8, 2014. ISSN 1472-6947. doi: 10.1186/1472-6947-14-S1-S3. URL <http://dx.doi.org/10.1186/1472-6947-14-S1-S3>.
- Fei Yu, Stephen E. Fienberg, Aleksandra B. Slavkovic, and Caroline Uhler. Scalable privacy-preserving data sharing methodology for genome-wide association studies. *Journal of Biomedical Informatics*, 50:133–141, 2014.
- Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society, Series B*, 67:301–320, 2005.